

	<p>SAEED SALEHI Department of Mathematics University of Tabriz P.O.Box 51666-17766 Tabriz, Iran</p>	<p>Tel: +98 (0)411 339 2905 Fax: +98 (0)411 334 2102 E-mail: /root@SaeedSalehi.ir/ /SalehiPour@TabrizU.ac.ir/ Web: http://SaeedSalehi.ir/</p>	<p>$\oint_{\Sigma_{\alpha\ell\epsilon\hbar}}^{\Sigma_{\alpha\epsilon\epsilon\partial}} \cdot \text{ir}$</p>
---	---	--	--

Herbrand Consistency of Some Finite Fragments of Bounded Arithmetical Theories

Abstract

We formalize the notion of Herbrand Consistency in an appropriate way for bounded arithmetics, and show the existence of a finite fragment of $\text{I}\Delta_0$ whose Herbrand Consistency is not provable in the theory $\text{I}\Delta_0$. We also show the existence of an $\text{I}\Delta_0$ –derivable Π_1 –sentence such that $\text{I}\Delta_0$ cannot prove its Herbrand Consistency.

Acknowledgements This research is partially supported by grant No 89030062 of the Institute for Research in Fundamental Sciences (IPM), Tehran, Iran.

2010 Mathematics Subject Classification: 03F40 · 03F25 · 03F30.

Keywords: Herbrand Consistency · Bounded Arithmetic · Gödel's Second Incompleteness Theorem.

Date: 09 October 2011 (09.10.11)

1 Introduction

A consequence of Gödel's Second Incompleteness Theorem is Π_1 –separation of some mathematical theories; for example ZFC is not Π_1 –conservative over PA since $ZFC \vdash \text{Con}(\text{PA})$ but (by Gödel's theorem) $PA \not\vdash \text{Con}(\text{PA})$, where Con is the consistency predicate. Inside PA, the hierarchy $\{\text{I}\Sigma_n\}_{n \geq 0}$ is not Π_1 –conservative, since $\text{I}\Sigma_{n+1} \vdash \text{Con}(\text{I}\Sigma_n)$ (but again $\text{I}\Sigma_n \not\vdash \text{Con}(\text{I}\Sigma_n)$). As for the bounded arithmetics, we only know that the elementary arithmetic $I\Delta_0 + \text{Exp}$ is not Π_1 –conservative over $I\Delta_0 + \bigwedge_j \Omega_j$ (see Corollary 5.34 of [5]). One candidate for Π_1 –separating $I\Delta_0 + \text{Exp}$ from $I\Delta_0$ was the Cut-Free Consistency of $I\Delta_0$ (see [7]): it was already known that $I\Delta_0 + \text{Exp} \vdash \text{CFCCon}(I\Delta_0)$ and it was presumed that $I\Delta_0 \not\vdash \text{CFCCon}(I\Delta_0)$, where CFCCon stands for Cut-Free Consistency. Though this presumption took rather a long time to be established (see [14]), it opened a new line of research.

The problem of provability (or unprovability) of the cut-free consistency of weak arithmetics is an interesting (double) generalization of Gödel's Second Incompleteness Theorem: the theory (being restricted to bounded or weak arithmetics) and also the consistency predicate are both weakened. Here, we do not intend to outline the history of this research line, and refer the reader to [11, 12]. Nevertheless, we list some prominent results obtained so far, to put our new result in perspective.

Herbrand Consistency is denoted by HCon and (Semantic) Tableau Consistency by TabCon . Adamowicz (with Zbierski in 2001 [2] and) in 2002 [3] showed that $I\Delta_0 + \Omega_m \not\vdash \text{HCon}(I\Delta_0 + \Omega_m)$ for $m \geq 2$. She had already shown the unprovability $I\Delta_0 + \Omega_1 \not\vdash \text{TabCon}(I\Delta_0 + \Omega_1)$ in 1996 (but appeared in 2001 as [1]). Salehi improved the result of [3] in [10] by showing that $I\Delta_0 + \Omega_1 \not\vdash \text{HCon}(I\Delta_0 + \Omega_1)$ (see also [12]) and the result of [2] in [9, 10] by showing $S \not\vdash \text{HCon}(S)$ where S is an $I\Delta_0$ –derivable Π_2 –sentence. This result also implied that $I\Delta_0 \not\vdash \text{HCon}(\overline{I\Delta_0})$ holds for a re-axiomatization $\overline{I\Delta_0}$ of $I\Delta_0$. Willard [13] showed in 2002 that $I\Delta_0 \not\vdash \text{TabCon}(I\Delta_0)$ and also $I\Delta_0 \not\vdash \text{HCon}(I\Delta_0 + \Omega_0)$, where Ω_0 is the axiom of the totality of the squaring function $\Omega_0 : \forall x \exists y [y = x \cdot x]$. This was improved in [12] by showing $I\Delta_0 \not\vdash \text{HCon}(I\Delta_0)$, without using the Ω_0 axiom. It was also proved in [13] that $V \not\vdash \text{HCon}(V)$ for an $I\Delta_0$ –derivable Π_1 –sentence V . Kołodziejczyk [6] showed in 2006 that the unprovability $I\Delta_0 + \bigwedge_j \Omega_j \not\vdash \text{HCon}(I\Delta_0 + \Omega_1)$ holds; his result was stronger in a sense that it showed $I\Delta_0 + \bigwedge_j \Omega_j \not\vdash \text{HCon}(S + \Omega_1)$ for a finite fragment $S \subseteq I\Delta_0$.

In this paper we use an idea of an anonymous referee of [12] for defining evaluations in a more effective way (Definition 2.9) suitable for bounded arithmetics; this is a great step forward, noting our mentioning in [12] that “[o]ur definition of Herbrand Consistency is not best suited for $I\Delta_0$ ”. We then partially answer the question proposed by the anonymous referee of [11] (see Conjecture 4.1 in [11]). The author is grateful to both the referees, for suggestions and inspirations.

We show the existence of a finite fragment T of $I\Delta_0$ such that $I\Delta_0 \not\vdash \text{HCon}(T)$; this generalizes the result of [12]. We also show the existence of an $I\Delta_0$ –derivable Π_1 –sentence U such that $I\Delta_0 \not\vdash \text{HCon}(U)$; this generalizes the main result of [9, 10] and [13]. For keeping the paper short, and to avoid repeating some technical details, we apologetically invite the reader to consult [11, 12]. We also assume familiarity with the Bible of this field [5].

2 Herbrand Consistency of Arithmetical Theories

For getting a unique Skolemized formula, it is more convenient to negation normalize and rectify it.

Definition 2.1 (Rectified Negation Normal Form) *A formula is in negation normal form when no implication symbol \rightarrow appears in it, and the negation symbol \neg appears behind the atomic formulas only. A formula is rectified when different quantifiers refer to different variables and no variable appears both free and bound in the formula.* \diamond

Any formula can be uniquely negation normalized by removing the implication connectives (replacing formulas of the form $A \rightarrow B$ with $\neg A \vee B$) and then pushing the negations inside the sub-formulas by de Morgan's Law, until they get to the atomic formulas. Renaming the variables can rectify any formula. Thus one can negation normalize and rectify a formula uniquely, up to a variable renaming.

Definition 2.2 (Skolemization) For any existential formula $\exists x A(x)$ with $m(\geq 0)$ free variables, let $f_{\exists x A(x)}$ be a new m -ary function symbol (which does not occur in A ; cf. [4]). For any rectified negation normal formula φ we define φ^S inductively:

- $\varphi^S = \varphi$ for atomic or negated-atomic formula φ
- $(\varphi \wedge \psi)^S = \varphi^S \wedge \psi^S$
- $(\varphi \vee \psi)^S = \varphi^S \vee \psi^S$
- $(\forall x \varphi)^S = \forall x \varphi^S$
- $(\exists x \varphi)^S = \varphi^S[f_{\exists x \varphi(x)}(\bar{y})/x]$ where \bar{y} are the free variables of $\exists x \varphi(x)$.

Finally, the Skolemized form φ^{Sk} of the formula φ is obtained by removing all the (universal) quantifiers of φ^S . The resulted formula is an open (quantifier-less) formula, with probably some free variables. If those (free) variables are substituted with some ground (variable-free) terms, we obtain an Skolem instance of that formula. \diamond

Summing up, to get an Skolem instance of a given formula φ we first negation normalize and then rectify it to get a formula φ^{RNNF} ; then we remove the quantifiers of $(\varphi^{\text{RNNF}})^S$ to get $(\varphi^{\text{RNNF}})^{\text{Sk}}$, and substituting its free variables with some ground terms, gives us an Skolem instance of the formula φ . Let us note that the Skolem instances of a formula are determined uniquely.

Theorem 2.3 (Herbrand-Skolem-Gödel) Any theory T is equi-consistent with its Skolemized theory. In other words, the theory T is consistent if and only if every finite set of Skolem instances of T is (propositionally) satisfiable. \square

Example 2.4 In the language of arithmetic $\mathcal{L}_A = \{0, S, +, \cdot, \leq\}$, let Ind_{\square} be the instance of induction principle $\psi(0) \wedge \forall x[\psi(x) \rightarrow \psi(S(x))] \rightarrow \forall x \psi(x)$ for $\psi(x) = \exists y[y \leq x \cdot x \wedge y = x \cdot x]$. This is an axiom of the theory $\text{I}\Delta_0$. Rectified Negation Normal Form (Ind_{\square})^{RNNF} of Ind_{\square} is

$$\forall u[u \not\leq 0 \cdot 0 \vee u \neq 0 \cdot 0] \vee \exists w \left[\exists z[z \leq w \cdot w \wedge z = w \cdot w] \wedge \forall v[v \not\leq S(w) \cdot S(w) \vee v \neq S(w) \cdot S(w)] \right] \vee \forall x \exists y[y \leq x \cdot x \wedge y = x \cdot x].$$

If \mathbf{c} is the Skolem constant symbol for $\exists w \left[\exists z[z \leq w \cdot w \wedge z = w \cdot w] \wedge \forall v[v \not\leq S(w) \cdot S(w) \vee v \neq S(w) \cdot S(w)] \right]$, and $\mathbf{q}(x)$ is the Skolem function symbol for the formula $\exists z[z \leq x \cdot x \wedge z = x \cdot x]$, then $((\text{Ind}_{\square})^{\text{RNNF}})^S$ is

$$\forall u[u \not\leq 0 \cdot 0 \vee u \neq 0 \cdot 0] \vee [[\mathbf{q}(\mathbf{c}) \leq \mathbf{c} \cdot \mathbf{c} \wedge \mathbf{q}(\mathbf{c}) = \mathbf{c} \cdot \mathbf{c}] \wedge \forall v[v \not\leq S(\mathbf{c}) \cdot S(\mathbf{c}) \vee v \neq S(\mathbf{c}) \cdot S(\mathbf{c})]] \vee \forall x[\mathbf{q}(x) \leq x \cdot x \wedge \mathbf{q}(x) = x \cdot x].$$

Finally, the Skolemized form $(\text{Ind}_{\square})^{\text{Sk}}$ of φ is obtained as:

$$[u \not\leq 0 \cdot 0 \vee u \neq 0 \cdot 0] \vee [[\mathbf{q}(\mathbf{c}) \leq \mathbf{c} \cdot \mathbf{c} \wedge \mathbf{q}(\mathbf{c}) = \mathbf{c} \cdot \mathbf{c}] \wedge [v \not\leq S(\mathbf{c}) \cdot S(\mathbf{c}) \vee v \neq S(\mathbf{c}) \cdot S(\mathbf{c})]] \vee [\mathbf{q}(x) \leq x \cdot x \wedge \mathbf{q}(x) = x \cdot x].$$

Substituting $u/0$, $v/S(\mathbf{c}) \cdot S(\mathbf{c})$, x/t will result in the following Skolem instance of φ :

$$[0 \not\leq 0 \cdot 0 \vee 0 \neq 0 \cdot 0] \vee [[\mathbf{q}(\mathbf{c}) \leq \mathbf{c} \cdot \mathbf{c} \wedge \mathbf{q}(\mathbf{c}) = \mathbf{c} \cdot \mathbf{c}] \wedge [S(\mathbf{c}) \cdot S(\mathbf{c}) \not\leq S(\mathbf{c}) \cdot S(\mathbf{c}) \vee S(\mathbf{c}) \cdot S(\mathbf{c}) \neq S(\mathbf{c}) \cdot S(\mathbf{c})]] \vee [\mathbf{q}(t) \leq t \cdot t \wedge \mathbf{q}(t) = t \cdot t]. \diamond$$

Propositional satisfiability is usually arithmetized from the usual provability, only in propositional logic (see e.g. [5]); but in a series of more recent papers, this notion have been arithmetized differently, according to ones needs ([1, 2, 3, 6, 9, 10, 11, 12, 13]). We formalize the notion of propositional satisfiability by means of evaluations (as in the op. cit. papers) on sets of (Skolem) ground terms, but in a more effective way. To get a small evaluation on a given set of terms, we first sort its members, and then require the equality relation to be a congruence.

We will call the ground terms constructed from Skolem function (and constant) symbols, simply *terms*. For a set Λ , its cardinality will be denoted by $|\Lambda|$, and for a sequence s , its length will be also denoted by $|s|$. The $(i+1)$ th member of s is denoted by $(s)_i$ for any $i < |s|$; so $s = \langle (s)_0, (s)_1, \dots, (s)_{|s|-1} \rangle$. Let \approx and \prec be two new symbols, not in the language of arithmetic $\mathcal{L}_A = \langle 0, S, +, \cdot, \leq \rangle$.

Definition 2.5 (Pre-Evaluation) For a set of terms Λ (with $|\Lambda| \geq 2$), a pre-evaluation on Λ is a sequence p that satisfies the following conditions:

- (1) length of p is $|p| = 2|\Lambda| - 1$;
- (2) for any $0 \leq i \leq |\Lambda| - 1$ we have $(p)_{2i} \in \Lambda$;
- (3) for any $1 \leq i \leq |\Lambda| - 1$ we have $(p)_{2i-1} \in \{\prec, \approx\}$;
- (4) for any term $t \in \Lambda$ there exists a unique $0 \leq j \leq |\Lambda| - 1$ such that $(p)_{2j} = t$. ◊

In other words, a pre-evaluation on Λ sorts (organizes) the terms in Λ , starting from the smallest and ending in the largest.

Example 2.6 A pre-evaluation on $\{\alpha_0, \alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5, \alpha_6\}$ is a sequence like

$$p = \langle \alpha_4, \prec, \alpha_7, \approx, \alpha_1, \approx, \alpha_5, \prec, \alpha_3, \prec, \alpha_6, \approx, \alpha_2 \rangle. \quad \diamond$$

Definition 2.7 (Equality and Order in Pre-Evaluations) In a pre-evaluation p on Λ define the relations \approx_p and \prec_p on Λ^2 by the following conditions for $s, t \in \Lambda$:

- (1) $s \approx_p t$ if there exists a sub-sequence q of p of length $2l - 1$ ($l \geq 1$) such that
 - (a) either $((q)_0 = s \& (q)_{2l-2} = t)$ or $((q)_0 = t \& (q)_{2l-2} = s)$;
 - (b) for any $1 \leq i \leq l - 1$, $(q)_{2i-1} = \approx$.
- (2) $s \prec_p t$ if there exists a sub-sequence q of p of length $2l - 1$ ($l \geq 1$) such that
 - (a) $(q)_0 = s$ and $(q)_{2l-2} = t$;
 - (b) there exists some $1 \leq i \leq l - 1$ for which $(q)_{2i-1} = \prec$. ◊

Example 2.6 (Continued) We have $\alpha_1 \approx_p \alpha_5 \approx_p \alpha_7$ and $\alpha_2 \approx_p \alpha_6$. Also, $\alpha_4 \prec_p \alpha_1$, $\alpha_4 \prec_p \alpha_5$, $\alpha_4 \prec_p \alpha_7$, $\alpha_1 \prec_p \alpha_2$, $\alpha_1 \prec_p \alpha_3$, and $\alpha_1 \prec_p \alpha_6$ hold. ◊

Lemma 2.8 (Equivalence and Order by Pre-Evaluation) Let Λ be a set of terms, and p be a pre-evaluation on Λ .

- (1) The relation \approx_p is an equivalence on Λ .
- (2) The relation \prec_p is a total order on Λ .
- (3) The relations \approx_p and \prec_p are compatible with each other: if $t \approx_p s$, and $t \prec_p u$ (respectively, $u \prec_p t$), then $s \prec_p u$ (respectively, $u \prec_p s$).

Proof. The parts (1) and (2) are immediate. For (3), suppose $t \approx_p s$ and $t \prec_p u$. Then there is a sub-sequence q of p which starts from t and ends with u and contains at least one special symbol \prec . There must also be some other sub-sequence r which starts from either t or s and ends with the other one, and all its special symbols are equality \approx . If r starts from s (and so ends with t), then the concatenation of r and q results in a sub-sequence which starts from s and ends with u and contains some special symbol \prec . Whence $s \prec_p u$. And if r starts from t , then q cannot be a sub-sequence of r because all the special symbols in r are \approx and q contains at least one special symbol \prec . Thus r has to be a sub-sequence of q . Then there must exist a sub-sequence of p which starts from s and ends with u and contains a special symbol \prec ; whence $s \prec_p u$. The other case ($u \prec_p t$) can be proved very similarly. \square

Definition 2.9 (Evaluation) A pre-evaluation p on a set of terms Λ is called an evaluation when, for any term $t, s \in \Lambda$ and any term $u(x)$ with the free variable x , if $t \approx_p s$ and $u(t/x), u(s/x) \in \Lambda$ hold, then $u(t/x) \approx_p u(s/x)$ holds too. \diamond

In other words, an evaluation on Λ is a pre-evaluation p on Λ whose equivalence relation \approx_p is a congruence relation on Λ .

Definition 2.10 (Satisfaction in an Evaluation) Let Λ be a set of terms and p an evaluation on it. For terms $t, s \in \Lambda$ we write $p \models t = s$ when $t \approx_p s$ holds. We also write $p \models t \leq s$ when either $t \approx_p s$ or $t \prec_p s$ holds. So, for atomic formulas φ in the language of arithmetic \mathcal{L}_A we have defined the notion of satisfaction $p \models \varphi$. The satisfaction relations can be extended to all open (quantifier-less) formulas as:

- $p \models \varphi \wedge \psi \iff p \models \varphi$ and $p \models \psi$
- $p \models \varphi \vee \psi \iff p \models \varphi$ or $p \models \psi$
- $p \models \varphi \rightarrow \psi \iff$ if $p \models \varphi$ then $p \models \psi$
- $p \models \neg \varphi \iff p \not\models \varphi$

\diamond

Lemma 2.11 (Leibniz's Law) Any evaluation p on any set of terms Λ satisfies all the available Skolem instances of the axioms of equational logic, in particular Leibniz's Law: for any $t, s \in \Lambda$ and any open formula $\varphi(x)$, we have $p \models t = s \wedge \varphi(t) \rightarrow \varphi(s)$.

Proof. Suppose $p \models t = s$. By induction on (the complexity) of (the open formula) φ one can show that $p \models \varphi(t)$ if and only if $p \models \varphi(s)$. For atomic φ it follows from Lemma 2.8, and for the more complex formulas it follows from the inductive definition of satisfaction in evaluations. \square

Definition 2.12 (T -evaluation on Λ) For a set of terms Λ , an Skolem instance of a formula is called to be available in Λ if all the terms appearing in it belong to Λ . For a theory T and a set of terms Λ and an evaluation p on Λ , we say that p is an T -evaluation on Λ if p satisfies every Skolem instance of every sentence in T which is available in Λ . \diamond

So, T -evaluations, for a theory T , are kind of partial models of T . Indeed, if Λ is the set of all (ground) terms (constructed from the language of T and the Skolem function symbols of the axioms of T), then any T -evaluaton on Γ (if exists) is a *Herbrand Model* of T . Herbrand's Theorem can be read as “A theory T is consistent if and only if for every finite set of (Skolem) terms, there exists an T -evaluation on it.” Thus, the notion of *Herbrand Consistency* of a theory T is (equivalent to) the existence of an T -evaluation on any (finite) set of terms.

Example 2.13 Let T be axiomatized by the following sentences in \mathcal{L}_A :

- $\forall x[x \cdot 0 = 0]$;
- $\exists y \leq 0 \cdot 0 [y = 0 \cdot 0] \wedge \forall x [\exists y \leq x \cdot x [y = x \cdot x] \rightarrow \exists y \leq S(x) \cdot S(x) [y = S(x) \cdot S(x)]] \rightarrow \forall x \exists y \leq x \cdot x [y = x \cdot x]$.

Let $\Lambda = \{0, 0 \cdot 0, \mathbf{c}, \mathbf{c} \cdot \mathbf{c}, \mathbf{q}(\mathbf{c}), S(\mathbf{c}) \cdot S(\mathbf{c}), t, t \cdot t, \mathbf{q}(t)\}$ where \mathbf{c} and \mathbf{q} are as in Example 2.4. As we saw in that example, the following is an instance of the the second axiom (Ind_\square), which is also available in Λ :

$$[0 \leq 0 \cdot 0 \vee 0 \neq 0 \cdot 0] \vee \\ [[\mathbf{q}(\mathbf{c}) \leq \mathbf{c} \cdot \mathbf{c} \wedge \mathbf{q}(\mathbf{c}) = \mathbf{c} \cdot \mathbf{c}] \wedge [S(\mathbf{c}) \cdot S(\mathbf{c}) \leq S(\mathbf{c}) \cdot S(\mathbf{c}) \vee S(\mathbf{c}) \cdot S(\mathbf{c}) \neq S(\mathbf{c}) \cdot S(\mathbf{c})]] \vee \\ [\mathbf{q}(t) \leq t \cdot t \wedge \mathbf{q}(t) = t \cdot t].$$

Suppose p is an T -evaluation on Λ . By the first axiom p must satisfy the instance $0 \cdot 0 = 0$, so we should have $p \models 0 \cdot 0 = 0$. Thus, p cannot satisfy the first disjunct of the above instance. Indeed, p cannot satisfy the second disjunct either, because for any term u we have $p \models u \leq u \wedge u = u$. Thus, p cannot satisfy the second conjunct of the second disjunct. Whence, p must satisfy the third disjunct of the above instance, and in particular we should have $p \models \mathbf{q}(t) = t \cdot t$. \diamond

Definition 2.14 (Skolem Hull) Let $\mathcal{L}_A^{\text{Sk}}$ be the language expanding \mathcal{L}_A by the Skolem function (and constant) symbols of all the existential formulas in the language \mathcal{L}_A . Or in other words, $\mathcal{L}_A^{\text{Sk}}$ is the set $\mathcal{L}_A^{\text{Sk}} = \{f_{\exists x \varphi(x)} \mid \varphi \text{ is an } \mathcal{L}_A\text{-formula}\}$. For a given set of terms Λ , let $\Lambda^{\langle j \rangle}$ be defined by induction on j :

$\Lambda^{\langle 0 \rangle} = \Lambda$, and

$\Lambda^{\langle j+1 \rangle} = \Lambda^{\langle j \rangle} \cup \{f(t_1, \dots, t_m) \mid f \in \mathcal{L} \wedge t_1, \dots, t_m \in \Lambda^{\langle j \rangle}\} \cup \{f_{\exists x \varphi(x)}(t_1, \dots, t_m) \mid \ulcorner \varphi \urcorner \leq j \wedge t_1, \dots, t_m \in \Lambda^{\langle j \rangle}\}$, where $\ulcorner \varphi \urcorner$ is the Gödel code of φ . \diamond

Bounding the Gödel code of φ in the above definition will enable us to have some efficient (upper bound) for the Gödel code of $\Lambda^{\langle j \rangle}$ (see [11, 12]).

Herbrand's theorem implies that for any \exists -formula $\exists x \psi(x)$ (where ψ is an open formula) and any theory T , if $T \vdash \exists x \psi(x)$ then there are some (Skolem) terms t_1, \dots, t_n such that $T^{\text{Sk}} \vdash \psi(t_1) \vee \dots \vee \psi(t_n)$. Usually this observation is called Herbrand's Theorem. We will need a somehow dual of this fact.

Lemma 2.15 (Herbrand Proof of Universal Formulas) For a \forall -formula $\forall x \psi(x)$ (where ψ is open) and a theory T , suppose $T \vdash \forall x \psi(x)$. Let Λ be a set of terms and $t \in \Lambda$. There exists a finite (standard) $k \geq 0$ such that for any T -evaluation p on $\Lambda^{\langle k \rangle}$ we have $p \models \psi(t)$.

Proof. By $T \vdash \forall x \psi(x)$ the theory $T^{\text{Sk}} \cup \{\neg \psi(\mathbf{c})\}$, where \mathbf{c} is the Skolem constant symbol for $\exists x \neg \psi(x)$, is inconsistent. Suppose φ is the rectified negation normal form of $\neg \psi$. Then, by Herbrand's theorem, there exists some finite set of terms Γ such that there can be no $(T^{\text{Sk}} \cup \{\varphi(\mathbf{c})\})$ -evaluation on it. Since \mathbf{c} appears in Γ we write it as $\Gamma(\mathbf{c})$, and by $\Gamma(t)$ we denote the set of terms which result from the terms of $\Gamma(\mathbf{c})$ by replacing \mathbf{c} with t everywhere. It can be clearly seen that there exists some $k \in \mathbb{N}$ such that $\Gamma(t) \subseteq \Lambda^{\langle k \rangle}$. Whence, there cannot be any $(T^{\text{Sk}} \cup \{\varphi(t)\})$ -evaluation on $\Lambda^{\langle k \rangle}$. Thus, any T -evaluation p on $\Lambda^{\langle k \rangle}$ must satisfy $p \not\models \varphi(t)$ or $p \models \psi(t)$. \square

Example 2.16 Let the theory T , in the language of arithmetic \mathcal{L}_A , be axiomatized by

- (1) $\forall x[S(x) \neq 0]$
- (2) $\forall x, y[x + S(y) = S(x + y)]$
- (3) $\forall x \exists z[x \neq 0 \rightarrow x = S(z)]$
- (4) $\forall x, y \exists z[x \leq y \rightarrow z + x = y]$

For the open formula $\psi(x) = (x \leq 0 \rightarrow x = 0)$ we have $T \vdash \forall x \psi(x)$.

Let $\mathbf{p}(x)$ be the Skolem function for the formula $\exists z[x = 0 \vee x = S(z)]$, and $\mathbf{h}(x, y)$ be the Skolem function for the formula $\exists z[x \leq y \vee z + x = y]$. Then the Skolemized form T^{Sk} of the theory T will be as:

$$\begin{array}{ll} (1') S(x) \neq 0 & (2') x + S(y) = S(x + y) \\ (3') x = 0 \vee x = S(\mathbf{p}(x)) & (4') x \leq y \vee \mathbf{h}(x, y) + x = y \end{array}$$

For a fixed term t let Γ_t be the following set of terms:

$$\Gamma_t = \{0, t, \mathbf{h}(t, 0), \mathbf{h}(t, 0) + t, \mathbf{p}(t), S(\mathbf{p}(t)), \mathbf{h}(t, 0) + \mathbf{p}(t), \mathbf{h}(t, 0) + S(\mathbf{p}(t)), S(\mathbf{h}(t, 0) + \mathbf{p}(t))\}.$$

Now we show that any T -evaluation p on Γ_t must satisfy $p \models \psi(t)$ or, equivalently, if $p \models t \leq 0$ then $p \models t = 0$. Assume $p \models t \leq 0$. Then by the fourth axiom we have $p \models \mathbf{h}(t, 0) + t = 0$. If $p \models t = 0$ does not hold, then $p \models t \neq 0$, so by the third axiom we have $p \models t = S(\mathbf{p}(t))$. Whence, $p \models \mathbf{h}(t, 0) + S(\mathbf{p}(t)) = 0$. On the other hand, by the second axiom, $p \models \mathbf{h}(t, 0) + S(\mathbf{p}(t)) = S(\mathbf{h}(t, 0) + \mathbf{p}(t))$. So, we infer that $p \models S(\mathbf{h}(t, 0) + \mathbf{p}(t)) = 0$, which is in contradiction with the first axiom. Thus, $p \models t = 0$ must hold, which shows that $p \models \psi(t)$. \diamond

As was mentioned before, for a consistent theory T there must exist some Herbrand Model of T .

Definition 2.17 (Definable Herbrand Models) Let Λ be a set of terms, and define its Skolem Hull to be $\Lambda^{\langle\infty\rangle} = \bigcup_{n \in \mathbb{N}} \Lambda^{\langle n \rangle}$ (see Definition 2.14). For an evaluation p on $\Lambda^{\langle\infty\rangle}$, let $\mathfrak{M}(\Lambda, p) = \{t/p \mid t \in \Lambda^{\langle\infty\rangle}\}$, where t/p is the equivalence class of the relation \approx_p containing t (cf. Lemma 2.8). Put the structure

$$\begin{array}{l} (1) f^{\mathfrak{M}(\Lambda, p)}(t_1/p, \dots, t_m/p) = f(t_1, \dots, t_m)/p, \\ (2) R^{\mathfrak{M}(\Lambda, p)} = \{(t_1/p, \dots, t_m/p) \mid p \models R(t_1, \dots, t_m)\}, \end{array}$$

on $\mathfrak{M}(\Lambda, p)$, for any m -ary function symbol f and any m -ary relation symbol R . \diamond

Lemma 2.18 (Herbrand Models by Evaluations) The structure on $\mathfrak{M}(\Lambda, p)$ is well-defined, and for a theory T , if p is an T -evaluation on Λ then $\mathfrak{M}(\Lambda, p) \models T$. \square

3 Bounded Arithmetic and Herbrand Consistency

By an efficient Gödel coding (see e.g. Chapter V of [5]) we can code sets, sequences (and so the syntactic concepts like Skolem function symbols, Skolem instances, evaluations, etc.) such that the following ([5]) hold for any sequences α, β :

- $|\alpha * \beta| \leq 64 \cdot (|\alpha| \cdot |\beta|)$, where $*$ denotes concatenation;
- $|\alpha| \leq \log(|\alpha|)$.

It follows that for any sets A, B we have $|A \cup B| \leq 64 \cdot (|A| \cdot |B|)$ and $|A| \leq \log(|A|)$. We write $X \in \mathcal{O}(Y)$ to indicate that $X \leq Y \cdot n + n$ for some $n \in \mathbb{N}$; that is X is linearly bounded by Y . The above (efficient) coding has the property that for any sequence $U = \langle u_1, \dots, u_l \rangle$ we have $\log(|U|) \in \mathcal{O}(\sum_i \log(|u_i|))$. For any evaluation p on a set of terms Λ it can be seen that $\log(|p|) \in \mathcal{O}(\log(|\Lambda|))$.

Let us note that all of the concepts introduced so far can be formalized in the language of arithmetic \mathcal{L}_A . Here we make the observation that, having an arithmetically definable set of terms Λ , the sets $\Lambda^{\langle j \rangle}$ are all definable in arithmetic (in terms of Λ and j), but the set $\Lambda^{\langle\infty\rangle}$ is not definable by an arithmetical formula. We will come to this point later. The arithmetical theory we are interested here is denoted by $I\Delta_0$ which is usually axiomatized by Robinson's arithmetic, in the language \mathcal{L}_A , plus the induction axiom for bounded formulas (see e.g. [5]).

In this section we prove our main result: the existence of a finite fragment $T \subseteq I\Delta_0$ whose Herbrand Consistency is not provable in $I\Delta_0$. As the exponential function $x \mapsto 2^x$ is not available (provably total) in $I\Delta_0$, then we denote by \log the set of elements x for which $\exp(x) = 2^x$ exists. Let us note that for a model \mathcal{M} , the set $\log(\mathcal{M})$ is the logarithm of the elements of \mathcal{M} . The set \log is closed under S and $+$, but not under \times , in $I\Delta_0$. We will use the term *cut* for any definable and downward closed set (not necessarily closed under S) in the arithmetical models. The formula “ $y = \exp(x)$ ” is expressible in \mathcal{L}_A , and $I\Delta_0$ can prove some of the basic properties of \exp (cf. [5]), though cannot prove its totality: $I\Delta_0 \not\vdash \forall x \exists y [y = \exp(x)]$. By \log^2 we denote the set of elements x for which $\exp^2(x) = 2^{2^x}$ exists; the superscripts on top of the functions denote the iteration. Similarly, $\log^n = \{x \mid \exists y [y = \exp^n(x)]\}$, where \exp^n denotes the n time iteration of the exponential function \exp .

We use a deep theorem in bounded arithmetic, which happens to be the very last theorem of [5]. It reads, in our terminology, as:

For any $k \geq 0$ there exists a bounded formula $\varphi(x)$ such that
 $I\Delta_0 + \Omega_1 \vdash \forall x \in \log^{k+1} \varphi(x)$, but $I\Delta_0 + \Omega_1 \not\vdash \forall x \in \log^k \varphi(x)$.

It can be clearly seen that the theorem also holds for $I\Delta_0$ instead of $I\Delta_0 + \Omega_1$, and for any cut I (and its logarithm $\log I = \{x \mid \exists y \in I [y = \exp(x)]\}$) instead of \log^k (and its logarithm \log^{k+1}); see also [3] and (Theorem 3.6 of) [11].

Theorem 3.1 (Π₁–Separation of Logarithmic Cuts) *For any cut I there exists a bounded formula $\varphi(x)$ such that $I\Delta_0 \cup \{\exists x \in I \varphi(x)\}$ is consistent, but $I\Delta_0 \cup \{\exists x \in \log I \varphi(x)\}$ is not consistent.* \square

We will find the desired finite fragment of $I\Delta_0$ (whose Herbrand Consistency is not provable in $I\Delta_0$) in three steps (the following subsections) before proving the main result (in the last subsection). For doing so, we will show that for sufficiently strong finite fragments of $I\Delta_0$, like T , if $I\Delta_0 \vdash \text{HCon}(T)$ then the consistency of the theory $I\Delta_0 \cup \{\exists x \in I \theta(x)\}$, for some suitable cut I and a suitable bounded formula θ , implies the consistency of the theory $T \cup \{\exists x \in \log I \theta(x)\}$. As we will see, this contradicts Theorem 3.1.

3.1 The First Finite Fragment

Assuming the consistency of the theory $I\Delta_0 \cup \{\exists x \in I \varphi(x), \text{HCon}(T)\}$, and inconsistency of the theory $T \cup \{\exists x \in \log I \varphi(x)\}$, we can construct a model \mathfrak{M} , from a given model $\mathcal{M} \models I\Delta_0 \cup \{\exists x \in I \varphi(x), \text{HCon}(T)\}$, such that $\mathfrak{M} \models T \cup \{\exists x \in \log I \varphi(x)\}$; which is in contradiction with the assumptions. For that, let us take a (hypothetical) model $\mathcal{M} \models I\Delta_0 \cup \{a \in I \wedge \varphi(a)\} \cup \{\text{HCon}(T)\}$ for some $a \in \mathcal{M}$. Then we form the set $\Gamma = \{\underline{0}, \underline{1}, \underline{2}, \dots, \omega_1(a)\}$ where \underline{i} is a term in \mathcal{L}_A representing the number i , defined inductively as $\underline{0} = 0$ and $\underline{i+1} = S(\underline{i})$. From the assumption $\mathcal{M} \models \text{HCon}(T)$ we find an T –evaluation p on $\Lambda^{(j)}$, for a suitable j and a suitable Λ which contains the above set Γ . Then we can form the model $\mathfrak{M}(\Lambda, p)$ and, by some technical details, show that $\mathfrak{M}(\Lambda, p) \models T + \exists x \in \log I \varphi(x)$. The bound $\omega_1(a)$ assures us that the set Γ contains the range of (the bounded) quantifiers in the (bounded) formula $\varphi(a)$. For the Gödel code of \underline{i} we have $\log(\Gamma \underline{i}) \in \mathcal{O}(\log(2^i))$ and so $\log(\Gamma \underline{i}) \in \mathcal{O}(\log(2^{(\omega_1(a))^2}))$ whence $\log(\Gamma \underline{i}) \in \mathcal{O}(\log(\exp^2(2(\log a)^2)))$. We need the closure of Γ under the Skolem function symbols of (a finite fragment of) $I\Delta_0$, that is $\Gamma^{(\infty)}$ (see Definitions 2.17 and 2.14). Since, unfortunately, that set is not definable, we consider the set $\Gamma^{(j)}$ for a non-standard j , which makes sense if $\Gamma \underline{i}$ (and so a) is non-standard. In case a is standard, then the proof becomes trivial (see below). For some non-standard j with $j \leq \log^4(\Gamma \underline{i})$ we can form the set $\Gamma^{(j)}$, in case $\omega_2(\Gamma \underline{i})$ exists (see [11, 12]). And finally we have $\log(\omega_2(\Gamma \underline{i})) \in \mathcal{O}(\log(\exp^2(4(\log a)^4)))$.

Definition 3.2 (The Cut \mathcal{I}) The cut \mathcal{I} is defined to be $\{x \mid \exists y[y = \exp^2(4(\log a)^4)]\}$, and its logarithm is $\log \mathcal{I} = \{x \mid \exists y[y = \exp^2(4a^4)]\}$. \diamond

Applying theorem 3.1 to the cut \mathcal{I} defined above, we find a (fixed) bounded formula θ and a finite fragment $T_0 \subseteq I\Delta_0$ such that the theory $I\Delta_0 \cup \{\exists x \in \mathcal{I}\theta(x)\}$ is consistent, but $T_0 \cup \{\exists x \in \log \mathcal{I}\theta(x)\}$ is not consistent.

Definition 3.3 (The First Fragment T_0) Let T_0 be a finite fragment of $I\Delta_0$ for which there exists a (fixed) bounded formula θ such that the theory $I\Delta_0 \cup \{\exists x \in \mathcal{I}\theta(x)\}$ is consistent, but $T_0 \cup \{\exists x \in \log \mathcal{I}\theta(x)\}$ is not consistent. Let \mathcal{M} be a (fixed) model such that $\mathcal{M} \models I\Delta_0 \cup \{\exists x \in \mathcal{I}\theta(x)\}$. \diamond

In the rest of the paper we will show that for a finite fragment T of $I\Delta_0$ extending T_0 we have that $\mathcal{M} \not\models \text{HCon}(T)$, where HCon is the predicate of Herbrand Consistency.

3.2 The Second Finite Fragment

The proof of the main result goes roughly as follows: if $\mathcal{M} \models \text{HCon}(T)$, for a finite fragment $T \subseteq I\Delta_0$ to be specified later, then there exists (in \mathcal{M}) some T -evaluation p on some $\Lambda^{(j)}$, where $\Lambda \supseteq \Gamma$ is to be specified later and Γ and j are as in the previous subsection. Whence we can form the model $\mathfrak{M}(\Lambda, p)$, for which we already have $\mathfrak{M}(\Lambda, p) \models T$. Our second finite fragment T_1 will have the property that if $T \supseteq T_1$ then $\mathfrak{M}(\Lambda, p) \models \theta_0(\underline{a}/p)$. The third finite fragment T_2 will have the property that if $T \supseteq T_2$ then we have $\mathfrak{M}(\Lambda, p) \models \underline{a}/p \in \log \mathcal{I}$. So, finally we will get the model $\mathfrak{M}(\Lambda, p)$ which satisfies $\mathfrak{M}(\Lambda, p) \models T + [\underline{a}/p \in \log \mathcal{I} \wedge \theta_0(\underline{a}/p)]$, or, in the other words, $\mathfrak{M}(\Lambda, p) \models T \cup \{\exists x \in \log \mathcal{I}\theta_0(x)\}$ which is in contradiction with (the choice of the first finite fragment) $T_0 \subseteq T$.

Definition 3.4 (The Second Fragment T_1) Let T_1 be a finite fragment of $I\Delta_0$ which can prove the following ($I\Delta_0$ -provable \forall^* -)sentences:

- $x + 0 = x$
- $x \cdot 0 = 0$
- $x \leq 0 \leftrightarrow x = 0$
- $x \leq y \vee y \leq x$
- $x \leq z + x$
- $x + z \leq y + z \rightarrow x \leq y$
- $x \neq y \leftrightarrow S(x) \leq y \vee S(y) \leq x$
- $x + S(y) = S(x + y)$
- $x \cdot S(y) = x \cdot y + x$
- $x \leq S(y) \leftrightarrow x = S(y) \vee x \leq y$
- $x \leq y \leq z \rightarrow x \leq z$
- $x \leq x + z$
- $z \neq 0 \wedge x \cdot z \leq y \cdot z \rightarrow x \leq y$
- $x \not\leq y \leftrightarrow S(y) \leq x$

and also can prove the following ($I\Delta_0$ -provable $\forall^* \exists^*$ -)sentences:

- $x \leq y \rightarrow \exists z[z + x = y]$
- $y \neq 0 \rightarrow \exists q, r[x = r + q \cdot y \wedge r \leq y]$

\diamond

Remark 3.5 It can be seen that T_1 can prove the following arithmetical sentences:

- $S(x) \neq 0$
- $S(x) = S(y) \rightarrow x = y$
- $S(x) \not\leq x$
- $x \neq 0 \rightarrow \exists y[x = S(y)]$

For a proof, first note that by $x \leq y \vee y \leq x$ we have $\forall u[u \leq u]$, and also from $x \leq z + x$ and $x + 0 = x$ we get $\forall u[0 \leq u]$. Now, if $S(u) = 0$, then $S(u) \leq 0$, and so by the axiom $x \not\leq y \leftrightarrow S(y) \leq x$ we get $0 \not\leq u$, contradiction! Also from the same axiom it follows that $u \not\leq u \leftrightarrow S(u) \leq u$, and thus $S(u) \not\leq u$. If $S(u) = S(v)$ and $u \neq v$ then by $x \neq y \leftrightarrow S(x) \leq y \vee S(y) \leq x$ we have either $S(u) \leq v$ or $S(v) \leq u$. If $S(u) \leq v$ then $S(v) \leq v$, contradiction! The other case is similar. Finally, assume $u \neq 0$. Then by

$x \leq 0 \leftrightarrow x = 0$ we have $u \leq 0$ and so the axiom $x \leq y \leftrightarrow S(y) \leq x$ implies that $S(0) \leq u$. Thus, by $x \leq y \rightarrow \exists z[z + x = y]$ we have $v + S(0) = u$ for some v . Then from $x + S(y) = S(x + y)$ and $x + 0 = x$ we conclude that $S(v) = u$. Q.E.D \diamond

The main property of T_1 is the following:

Theorem 3.6 (The Main Property of T_1) Suppose $\mathcal{M} \models \text{I}\Delta_0 + [a \in \mathcal{I} \wedge \theta(a)] + \text{HCon}(T)$ is a non-standard model where θ is a bounded formula and $a \in \mathcal{M}$ is non-standard and $T \vdash T_1$. If $p \in \mathcal{M}$ is an T -evaluation on $\Lambda^{\langle j \rangle}$ where Λ is a set of terms such that $\Lambda \supseteq \Gamma = \{\underline{i} \mid i \leq \omega_1(a)\}$ and j is a non-standard element of \mathcal{M} , then for any bounded formula $\varphi(x_1, \dots, x_n)$ and any elements $i_1, \dots, i_n \leq a$, $\mathcal{M} \models \varphi(i_1, \dots, i_n) \iff \mathfrak{M}(\Lambda, p) \models \varphi(\underline{i_1}/p, \dots, \underline{i_n}/p)$.

We prove the theorem by induction on (the complexity) of φ (see also [11, 12]).

Lemma 3.7 (Another Property of T_1) Suppose $\mathcal{K} \models T_1$ and $a \in \mathcal{K}$, and let t be a term in \mathcal{L}_A . For any $i_1, \dots, i_n \leq a$ in \mathcal{K} and any $b \in \mathcal{K}$, if $\mathcal{K} \models b \leq t(i_1, \dots, i_n)$ then there exists a term s and there are some $j_1, \dots, j_m \leq a$ such that $\mathcal{K} \models b = s(j_1, \dots, j_m)$.

Proof. By induction on t :

- $t = 0$: if $\mathcal{K} \models b \leq 0$ then by the T_1 -axiom $x \leq 0 \leftrightarrow x = 0$ we have $\mathcal{K} \models b = 0$.
- $t = S(t_1)$: if $\mathcal{K} \models b \leq S(t_1)$ then by $x \leq S(y) \leftrightarrow x = S(y) \vee x \leq y$ which is a T_1 -axiom, we have $\mathcal{K} \models b = S(t_1) \vee b \leq t_1$, and the result follows from the induction hypothesis.
- $t = t_1 + t_2$: if $\mathcal{K} \models b \leq t_1 + t_2$ then by the T_1 -axiom $x \leq y \vee y \leq x$ we have that $\mathcal{K} \models b \leq t_2 \vee t_2 \leq b$. If $\mathcal{K} \models b \leq t_2$ then the conclusion follows from the induction hypothesis. Otherwise if $\mathcal{K} \models t_2 \leq b$ then by $x \leq y \rightarrow \exists z[z + x = y]$ (another T_1 -axiom) there exists some $d \in \mathcal{K}$ such that $\mathcal{K} \models d + t_2 = b$. Thus $\mathcal{K} \models d + t_2 \leq t_1 + t_2$, whence by the T_1 -axiom $x + z \leq y + z \rightarrow x \leq y$ we have $\mathcal{K} \models d \leq t_1$, and the desired result follows from the induction hypothesis and the fact that $\mathcal{K} \models b = d + t_2$.
- $t = t_1 \cdot t_2$: assume $\mathcal{K} \models b \leq t_1 \cdot t_2$. If $\mathcal{K} \models t_2 = 0$ then $\mathcal{K} \models t_1 \cdot t_2 = 0$ by the T_1 -axiom $x \cdot 0 = 0$. And so $\mathcal{K} \models b \leq 0$ is reduced to the first case above. Now suppose $\mathcal{K} \models t_2 \neq 0$. Then by the T_1 -axiom $y \neq 0 \rightarrow \exists q, r[x = r + q \cdot y \wedge r \leq y]$ we have $\mathcal{K} \models b = r + q \cdot t_2 \wedge r \leq t_2$ for some $q, r \in \mathcal{K}$. By the T_1 -axiom $x \leq z + x$ we have $\mathcal{K} \models q \cdot t_2 \leq r + q \cdot t_2 = b \leq t_1 \cdot t_2$ and so from the T_1 -axiom $x \leq y \leq z \rightarrow x \leq z$ it follows that $\mathcal{K} \models q \cdot t_2 \leq t_1 \cdot t_2$, and the T_1 -axiom $z \neq 0 \wedge x \cdot z \leq y \cdot z \rightarrow x \leq y$ implies that $\mathcal{K} \models q \leq t_1$ (since $\mathcal{K} \models t_2 \neq 0$). Now, the desired conclusion follows from the induction hypothesis and $\mathcal{K} \models b = r + q \cdot t_2 \wedge r \leq t_2 \wedge q \leq t_1$. \square

Lemma 3.8 (Preservation of Atomic Formulas) With the assumptions of Theorem 3.6 for any atomic formula $\varphi(x_1, \dots, x_n)$ and any $i_1, \dots, i_n \leq a$, we have that

$$\mathcal{M} \models \varphi(i_1, \dots, i_n) \iff \mathfrak{M}(\Lambda, p) \models \varphi(\underline{i_1}/p, \dots, \underline{i_n}/p).$$

Proof. By the T_1 -axioms $x \neq y \leftrightarrow S(x) \leq y \vee S(y) \leq x$ and $x \leq y \leftrightarrow S(y) \leq x$ it suffices to prove the one direction only: $\mathcal{M} \models \varphi(i_1, \dots, i_n) \implies \mathfrak{M}(\Lambda, p) \models \varphi(\underline{i_1}/p, \dots, \underline{i_n}/p)$. If $\varphi = "t \leq s"$ for some \mathcal{L}_A -terms t and s , then $\mathcal{M} \models t \leq s$ implies the existence of some $b \in \mathcal{M}$ such that $\mathcal{M} \models b + t = s$. By the T_1 -axiom $x \leq x + z$, $\mathcal{M} \models b \leq s$ so by Lemma 3.7 there exists an \mathcal{L}_A -term r (and some $j_1, \dots, j_m \leq a$) such that $\mathcal{K} \models b = r$. Whence, $\mathcal{M} \models r + t = s$. So, noting that $\mathcal{M}, \mathfrak{M}(\Lambda, p) \models T_1$, it suffices to prove the lemma for the atomic formula φ of the form $\varphi = "t = s"$.

For that we first note that if $i_1, \dots, i_n \leq a$ then $t(i_1, \dots, i_n), s(i_1, \dots, i_n) \leq \omega_1(a)$ holds. Suppose we have $\mathcal{M} \models t(i_1, \dots, i_n) = s(i_1, \dots, i_n) = i$. We show by induction on (the complexity of) t that the

condition $\mathcal{M} \models t(i_1, \dots, i_n) = i$ implies $\mathfrak{M}(\Lambda, p) \models t(\underline{i}_1/p, \dots, \underline{i}_n/p) = \underline{i}/p$. Let us note that the statement $\mathfrak{M}(\Lambda, p) \models t(\underline{i}_1/p, \dots, \underline{i}_n/p) = \underline{i}/p$ is equivalent to $\mathcal{M} \models "p \models t(\underline{i}_1, \dots, \underline{i}_n) = \underline{i}"$. So, it suffices to show the equivalence $\mathcal{M} \models t(i_1, \dots, i_n) = i \leftrightarrow "p \models t(\underline{i}_1, \dots, \underline{i}_n) = \underline{i}"$ by induction on t . For $t = 0$ and $t = S(t_1)$ the result follows from the definition $\underline{0} = 0$ and $\underline{j+1} = S(\underline{j})$. And for $t = t_1 + t_2$ and $t = t_1 \cdot t_2$ the result follows from the T_1 -axioms $x + 0 = x$, $x + S(y) = S(x + y)$, $x \cdot 0 = 0$, and $x \cdot S(y) = x \cdot y + x$. \square

Hence, the lemma also holds for open formulas φ as well. For bounded formulas we note that the range of quantifiers of $\varphi(i_1, \dots, i_n)$ for $i_1, \dots, i_n \leq a$ is contained in the set $\{j \mid j \leq \omega_1(a)\}$. This is formally expressed in the following lemma.

Lemma 3.9 (End-Extension Property) *With the assumptions of Theorem 3.6, if for some $i \leq a$ and some term t we have $(\mathcal{M} \models) p \models t \leq \underline{i}$ then there exists some $j \leq i$ such that $(\mathcal{M} \models) p \models t = \underline{j}$.*

Proof. By induction on the term \underline{i} . For $i = 0$, if $p \models t \leq 0$ then by Lemma 2.15, and the T_1 -axiom $x \leq 0 \leftrightarrow x = 0$, we have $p \models t = 0 = \underline{0}$. For $i = S(j)$, if $p \models t \leq S(j)$ then by Lemma 2.15, and the T_1 -axiom $x \leq S(y) \leftrightarrow x = S(y) \vee x \leq y$, we must have that $p \models t = \underline{S(j)} \vee t \leq \underline{j}$. Now the conclusion follows from the induction hypothesis. \square

Now we can prove Theorem 3.6.

Proof. (of Theorem 3.6) By induction on (the complexity of the bounded formula) φ . As the lemma has been proved for open formulas φ , it suffices to show that if the lemma holds for the (bounded) formula φ then it also holds for the (bounded) formula $\exists x \leq t(i_1, \dots, i_n) \varphi(x, i_1, \dots, i_n)$ where t is an \mathcal{L}_A -term; in the other words:

- $\mathcal{M} \models \exists x \leq t(i_1, \dots, i_n) \varphi(x, i_1, \dots, i_n) \iff \mathfrak{M}(\Lambda, p) \models \exists x \leq t(\underline{i}_1/p, \dots, \underline{i}_n/p) \varphi(\underline{i}_1/p, \dots, \underline{i}_n/p)$.
- If $\mathcal{M} \models b \leq t(i_1, \dots, i_n) \wedge \varphi(b, i_1, \dots, i_n)$, for some $b \in \mathcal{M}$, then by Lemma 3.7 there are terms s and elements $j_1, \dots, j_m \leq a$ such that $\mathcal{M} \models b = s(j_1, \dots, j_m)$. So, we have $\mathcal{M} \models \varphi(s(j_1, \dots, j_m), i_1, \dots, i_n)$. Whence, by the induction hypothesis we also have $\mathfrak{M}(\Lambda, p) \models \varphi(s(\underline{j}_1/p, \dots, \underline{j}_m/p), \underline{i}_1/p, \dots, \underline{i}_n/p)$, thus, noting that we already have $\mathfrak{M}(\Lambda, p) \models s(\underline{j}_1/p, \dots, \underline{j}_m/p) \leq t(\underline{i}_1/p, \dots, \underline{i}_n/p)$, the desired conclusion holds: $\mathfrak{M}(\Lambda, p) \models \exists x \leq t(\underline{i}_1/p, \dots, \underline{i}_n/p) \varphi(\underline{i}_1/p, \dots, \underline{i}_n/p)$.
- Conversely, if $\mathfrak{M}(\Lambda, p) \models d \leq t(\underline{i}_1/p, \dots, \underline{i}_n/p) \wedge \varphi(d, \underline{i}_1/p, \dots, \underline{i}_n/p)$ holds for some $d \in \mathfrak{M}(\Lambda, p)$ then by Lemma 3.7 there are some \mathcal{L}_A -term s and some $l_1, \dots, l_m \leq \underline{a}/p$ such that $\mathfrak{M}(\Lambda, p) \models d = s(l_1, \dots, l_m)$. For each $\alpha \leq m$ there is some term $\ell_\alpha \in \Lambda^{(\infty)}$ such that $l_\alpha = \ell_\alpha/p$. For each such α we also have that $\mathfrak{M}(\Lambda, p) \models \ell_\alpha/p \leq \underline{a}/p$ or equivalently $\mathcal{M} \models "p \models \ell_\alpha \leq \underline{a}"$. So, by Lemma 3.9 there exists some $j_\alpha \leq a$ for which we have $\mathcal{M} \models \ell_\alpha = \underline{j}_\alpha$. Whence, $\mathfrak{M}(\Lambda, p) \models d = s(j_1/p, \dots, j_m/p)$ and so

$$\mathfrak{M}(\Lambda, p) \models s(j_1/p, \dots, j_m/p) \leq t(\underline{i}_1/p, \dots, \underline{i}_n/p), \text{ and}$$

$$\mathfrak{M}(\Lambda, p) \models \varphi(s(j_1/p, \dots, j_m/p), \underline{i}_1/p, \dots, \underline{i}_n/p).$$

Thus, by the induction hypothesis we have

$$\mathcal{M} \models s(j_1, \dots, j_m) \leq t(i_1, \dots, i_n), \text{ and } \mathcal{M} \models \varphi(s(j_1, \dots, j_m), i_1, \dots, i_n).$$

So, we conclude that $\mathcal{M} \models \exists x \leq t(i_1, \dots, i_n) \varphi(x, i_1, \dots, i_n)$. \square

Let us repeat where we are now: in looking for a finite fragment $T \subseteq I\Delta_0$ such that $I\Delta_0 \not\models \text{HCon}(T)$ we found a finite fragment $T_0 \subseteq I\Delta_0$ and a bounded formula $\theta(x)$ such that $T_0 \vdash \neg \exists x \in \log \mathcal{I} \theta(x)$ but the theory $I\Delta_0 + \exists x \in \mathcal{I} \theta(x)$ is consistent and has a model $\mathcal{M} \models I\Delta_0 + [a \in \mathcal{I} \wedge \theta(a)]$. Then we aim at showing that $\mathcal{M} \not\models \text{HCon}(T)$. If $\mathcal{M} \models \text{HCon}(T)$ then we form the set of formulas $\Gamma = \{\underline{i} \mid i \leq \omega_1(a)\}$ for which $\omega_2(\Gamma)$ exists (by the very definition of \mathcal{I} and the assumption $a \in \mathcal{I}$), and so we can form the model $\mathfrak{M}(\Gamma, p)$ where p is an T -evaluaiton on $\Gamma^{(j)}$ (where $j \leq \log^4(\Gamma)$ can be taken to be non-standard if a is

so). The theory T_1 had the property that $\mathfrak{M}(\Gamma, p) \models \theta(a/p)$ (by Theorem 3.6), and in the next subsection we introduce a finite fragment $T_2 \subseteq \text{I}\Delta_0$ such that for a suitable $\Lambda \supseteq \Gamma$ (to be defined later) we will have $\mathfrak{M}(\Lambda, p) \models \underline{a}/p \in \log \mathcal{I}$. Then by taking T to be any finite fragment of $\text{I}\Delta_0$ which extends $T_0 \cup T_1 \cup T_2$ we will conclude that $\mathcal{M} \models \neg \text{HCon}(T)$.

3.3 The Third Finite Fragment

The fragments T_0 and T_1 were chosen not by their axioms but by their implications; T_0 had to prove $\neg \exists x \in \log \mathcal{I} \theta(x)$ (Definition 3.3), and T_1 had to prove some certain arithmetical statements (Definition 3.4). But for T_2 we require that it contains one of the following sentences as (one of) its (explicit) axioms (not only its consequences).

Definition 3.10 (Axioms for Totality of Squaring Function) (1) *The induction principle for the bounded formula $\psi(x) = \exists y \leq x^2 [y = x \cdot x]$ is denoted by $\text{Ind}_\square : \psi(0) \wedge \forall x (\psi(x) \rightarrow \psi(S(x))) \rightarrow \forall x \psi(x)$. Or, in other words (cf. Examples 2.4, 2.13) Ind_\square , which is an axiom of the theory $\text{I}\Delta_0$, is the sentence:*

$$\exists y \leq 0^2 [y = 0 \cdot 0] \wedge \forall x (\exists y \leq x^2 [y = x \cdot x] \rightarrow \exists y \leq S(x)^2 [y = S(x) \cdot S(x)]) \implies \forall x \exists y \leq x^2 [y = x \cdot x].$$

(2) *The Π_1 -sentence expressing the totality of squaring is denoted by $\Omega_0 : \forall x \exists y \leq x^2 [y = x \cdot x]$. \diamond*

We denote by $\mathbf{q}(x)$ the Skolem function symbol of the formula $\exists y \leq x^2 [y = x \cdot x]$ (cf. Examples 2.4, 2.13). Then the Skolemized forms of the axioms of Definition 3.10 will be as

1. $[u \not\leq 0^2 \vee u \neq 0 \cdot 0] \vee$
 $[[\mathbf{q}(\mathbf{c}) \leq \mathbf{c}^2 \wedge \mathbf{q}(\mathbf{c}) = \mathbf{c} \cdot \mathbf{c}] \wedge [v \not\leq S(\mathbf{c})^2 \vee v \neq S(\mathbf{c}) \cdot S(\mathbf{c})]] \vee$
 $[\mathbf{q}(x) \leq x^2 \wedge \mathbf{q}(x) = x \cdot x],$

where u, v, x are free variables and \mathbf{c} is the Skolem constant as in Example 2.4.

2. $\mathbf{q}(x) \leq x^2 \wedge \mathbf{q}(x) = x \cdot x.$

Define the terms \mathbf{q}_i 's by induction: $\mathbf{q}_0 = S(S(0))$ and $\mathbf{q}_{i+1} = \mathbf{q}(\mathbf{q}_i)$. It can be easily seen that \mathbf{q}_i represents the number $\exp^2(i)$, while for the code of \mathbf{q}_i we have $\log(\lceil \mathbf{q}_i \rceil) \in \mathcal{O}(\log(\exp(i)))$. That is to say that while the value of the term \mathbf{q}_i is of double exponential, the code of it is of (single) exponential. This (one) exponential gap, will make our proof to go through.

Formulating the statement " $x \in \log^2$ " can be stated as "there exists a sequence s such that $(s)_0 = 2$ and $|s| = x+1$ and for any $i < x$ we have $(s)_{i+1} = (s)_i \cdot (s)_i$ ". And " $y \in \log \mathcal{I}$ " can be stated as " $4y^4 \in \log^2$ ". Put $\Upsilon = \{\mathbf{q}_i \mid i \leq 4a^4\}$. Then any Ω_0 -evaluaton or Ind_\square -evaluation on $\Upsilon^{(\infty)}$ must satisfy $\mathbf{q}_{i+1} = \mathbf{q}_i \cdot \mathbf{q}_i$ for any $i < 4a^4$. If p is any such evaluation, then $\mathfrak{M}(\Upsilon, p) \models \forall i < 4(a/p)^4 [\mathbf{q}_{i+1}/p = \mathbf{q}_i/p \cdot \mathbf{q}_i/p]$. We require the finite fragment $T_2 \subseteq \text{I}\Delta_0$ to have the property that for any model $\mathcal{K} \models T_2$ if there are elements $q_0, q_1, \dots, q_b \in \mathcal{K}$ such that \mathcal{K} satisfies $q_0 = 2$ and $q_{i+1} = q_i^2$ for any $i < b$, then $\mathcal{K} \models b \in \log^2$. Let us note that the code of the sequence $\langle \exp^2(0), \exp^2(1), \dots, \exp^2(b) \rangle$ is roughly bounded by $\prod_{i \leq b} \exp^2(i) \approx (\exp^2(b))^2 = \exp^2(b+1)$. So, in the presence of $q_0, q_1, \dots, q_b \in \mathcal{K}$ with the above property, the (code of the) sequence s with the property " $(s)_0 = 2$, $|s| = x+1$ and for any $i < x$, $(s)_{i+1} = (s)_i \cdot (s)_i$ " must exist. Note also that $\text{I}\Delta_0 \vdash \forall i [i \in \log^2 \rightarrow i+1 \in \log^2]$. $(*)$

Definition 3.11 (The Third Fragment T_2) (1) *If the usual axiomatization of $\text{I}\Delta_0$ is taken into account, then let T_2 be a finite fragment of it which contains the axiom Ind_\square and has the property $(*)$ above.*

(2) If $I\Delta_0$ has been axiomatized all by Π_1 -formulas, where the induction axioms are in the form

$$\forall y(\varphi(0) \wedge \forall x < y[\varphi(x) \rightarrow \varphi(S(x))] \rightarrow \forall x \leq y \varphi(x))$$

for bounded φ , then we take the theory T_2 to be a finite fragment of $I\Delta_0^\Pi + \Omega_0$, where $I\Delta_0^\Pi$ is the above Π_1 -axiomatization of $I\Delta_0$, together with the axiom Ω_0 , such that it has the property $(*)$ above. So, in this case T_2 is a Π_1 -theory. \diamond

Let us reiterate the main property of T_2 again.

The Main Property of T_2 For a model $\mathcal{K} \models T_2$ if there are $q_0, q_1, \dots, q_b \in \mathcal{K}$ such that for any $j < b$ we have $\mathcal{K} \models q_{j+1} = q_j^2$ then $\mathcal{K} \models "b \in \log^2"$. \diamond

3.4 The Proof of the Main Result

Let T be any finite fragment of $I\Delta_0$ or $I\Delta_0^\Pi + \Omega_0$ such that $T \supseteq T_0 \cup T_1 \cup T_2$. If T_2 is taken as in the clause (1) of Definition 3.11 then T is truly a finite fragment of $I\Delta_0$, and if T_2 is taken as in the clause (2) of Definition 3.11 then T is a finite $I\Delta_0$ -derivable Π_1 -theory, whose conjunction (denoted by U) is a $I\Delta_0$ -derivable Π_1 -sentence.

Theorem 3.12 (The Main Theorem) (1) For a finite fragment T of $I\Delta_0$ we have $I\Delta_0 \not\models \text{HCon}(T)$.

(2) There exists an $I\Delta_0$ -derivable Π_1 -sentence U such that $I\Delta_0 \not\models \text{HCon}(U)$.

Proof. For the part (1) take T_2 as in clause (1) of Definition 3.11, and for part (2) take T_2 as in clause (2) of Definition 3.11, and let U be the conjunction of the axioms of T . In each case we will have the Skolem function symbol $\mathbf{q}(x)$ for squaring $x \mapsto x^2$.

By Theorem 3.1 there exists a (fixed) bounded formula $\theta(x)$, for the cut \mathcal{I} defined in Definition 3.2, such that $I\Delta_0 \not\models \neg \exists x \in \mathcal{I} \theta(x)$ and $T_0 \vdash \neg \exists x \in \log \mathcal{I} \theta(x)$ (see Definition 3.3). Fix $\mathcal{M} \models I\Delta_0 + [a \in \mathcal{I} \wedge \theta(a)]$. We show that $\mathcal{M} \not\models \text{HCon}(T)$.

Assume, for the sake of contradiction, that $\mathcal{M} \models \text{HCon}(T)$. Define the terms \underline{i} 's and \mathbf{q}_i 's by induction: $\underline{0} = 0$, $\underline{i+1} = S(\underline{i})$, $\mathbf{q}_0 = \underline{2}$, $\mathbf{q}_{i+1} = \mathbf{q}(\mathbf{q}_i)$. Let Λ be the set of terms $\{\underline{i} \mid i \leq \omega_a(a)\} \cup \{\mathbf{q}_i \mid i \leq \omega_1(a)\}$ in \mathcal{M} . As we saw earlier, the code of \underline{i} (and \mathbf{q}_i) are bounded by some polynomial of $\exp(i)$ and the code of the Λ is polynomially bounded by $\exp((\omega_1(a)^2))$ or $\exp^2(2(\log a)^2)$, and finally $\omega_2(\Gamma \Lambda^\neg)$ is polynomially bounded by $\exp^2(4(\log a)^4)$; which exists by the assumption $a \in \mathcal{I}$. We note that a is non-standard, because otherwise we would have $a \in \log \mathcal{I}$ and whence \mathcal{M} would be a model of the inconsistent theory $I\Delta_0 + \exists x \in \log \mathcal{I} \theta(x)$; a contradiction with the hypothesis. The existence of $\omega_2(\Gamma \Lambda^\neg)$ assures the existence of a non-standard element $j (\leq \log^4(\Gamma \Lambda^\neg))$ for which $\Lambda^{(j)}$ exists, and so by the assumption $\mathcal{M} \models \text{HCon}(T)$ there must exist some T -evaluation p on $\Lambda^{(j)}$ (hence, on $\Lambda^{(\infty)}$) in \mathcal{M} . So, we can form the model $\mathfrak{M}(\Lambda, p)$. For this model we have $\mathfrak{M}(\Lambda, p) \models T$ by Lemma 2.18. Since $\mathcal{M} \models \theta(a)$ (and $\mathfrak{M}(\Lambda, p) \models T_1$) then $\mathfrak{M}(\Lambda, p) \models \theta(\underline{a}/p)$ by Theorem 3.6. Also, since $\mathfrak{M}(\Lambda, p) \models T_2$ and $\mathbf{q}_0, \mathbf{q}_1, \dots, \mathbf{q}_b$ (for $b = 4a^4$) are elements of $\mathfrak{M}(\Lambda, p)$ such that $\mathfrak{M}(\Lambda, p) \models \mathbf{q}_0 = \underline{2}$ and $\mathfrak{M}(\Lambda, p) \models \mathbf{q}_{i+1} = \mathbf{q}_i^2$ for any $i < b$, then (by the main property of T_2) $\mathfrak{M}(\Lambda, p) \models "b \in \log^2"$. Or, in other words, $\mathfrak{M}(\Lambda, p) \models "\underline{a}/p \in \log \mathcal{I}"$. Whence, $\mathfrak{M}(\Lambda, p) \models [\underline{a}/p \in \log \mathcal{I} \wedge \theta(\underline{a}/p)]$. So, $\mathfrak{M}(\Lambda, p)$ is a model of $T + \exists x \in \log \mathcal{I} \theta(x)$, and this is contradiction with the assumption of $T \supseteq T_0$ and the inconsistency of the theory $T_0 + \exists x \in \log \mathcal{I} \theta(x)$. Thus $\mathcal{M} \not\models \text{HCon}(T)$ and so $I\Delta_0 \not\models \text{HCon}(T)$. \square

References

- [1] ZOFIA ADAMOWICZ, On Tableaux Consistency in Weak Theories, Preprint # 618, Institute of Mathematics, Polish Academy of Sciences (2001) <http://www.impan.pl/Preprints/p618.ps>
- [2] ZOFIA ADAMOWICZ & PAWEŁ ZBIERSKI, On Herbrand Consistency in Weak Arithmetic, *Archive for Mathematical Logic* **40**, 399–413 (2001) <http://dx.doi.org/10.1007/s001530000072>
- [3] ZOFIA ADAMOWICZ, Herbrand Consistency and Bounded Arithmetic, *Fundamenta Mathematicae* **171**, 279–292 (2002) <http://journals.impan.gov.pl/fm/Inf/171-3-7.html>
- [4] SAMUEL R. BUSS, On Herbrand’s Theorem, in: Maurice, D., Leivant, R. (eds.): Selected Papers from the International Workshop on *Logic and Computational Complexity*, Indianapolis, IN, USA, October 13–16, 1994, Lecture Notes in Computer Science, vol. 960, Springer-Verlag (1995) pp. 195–209
<http://math.ucsd.edu/~sbuss/ResearchWeb/herbrandtheorem/>
- [5] PETR HÁJEK & PAVEL PUDLÁK, *Metamathematics of First-Order Arithmetic*, Springer-Verlag, 2nd printing (1998) <http://projecteuclid.org/handle/euclid.pl/1235421926>
- [6] LESZEK ALEKSANDER KOŁODZIEJCZYK, On the Herbrand Notion of Consistency for Finitely Axiomatizable Fragments of Bounded Arithmetic Theories, *Journal of Symbolic Logic* **71**, 624–638 (2006)
<http://dx.doi.org/10.2178/jsl/1146620163>
- [7] JEFF B. PARIS & ALEX J. WILKIE, Δ_0 Sets and Induction, in: Guzicki W. & Marek W. & Plec A. & Rauszer C. (eds.) *Proceedings of Open Days in Model Theory and Set Theory*, Jadwisin, Poland 1981, Leeds University Press (1981) pp. 237–248
- [8] PAVEL PUDLÁK, Cuts, Consistency Statements and Interpretations, *Journal of Symbolic Logic* **50**, 423–441 (1985) <http://www.jstor.org/stable/2274231>
- [9] SAEED SALEHI, Unprovability of Herbrand Consistency in Weak Arithmetics, in: Striegnitz K. (ed.), Proceedings of the sixth ESSLLI Student Session, *European Summer School for Logic, Language, and Information* (2001) pp. 265–274 <http://saeedsalehi.ir/pdf/esslli.pdf>
- [10] SAEED SALEHI, *Herbrand Consistency in Arithmetics with Bounded Induction*, Ph.D. Dissertation, Institute of Mathematics, Polish Academy of Sciences (2002) <http://saeedsalehi.ir/pphd.html>
- [11] SAEED SALEHI, Separating bounded arithmetical theories by Herbrand consistency, *Journal of Logic and Computation* (to appear) <http://dx.doi.org/10.1093/logcom/exr005>
Preprint arXiv:1008.0225v2 [math.LO] (2010) <http://arxiv.org/pdf/1008.0225v2>
- [12] SAEED SALEHI, Herbrand Consistency of Some Arithmetical Theories, Submitted for publication.
Preprint arXiv:1005.2654v2[math.LO] (2010) <http://arxiv.org/pdf/1005.2654>
- [13] DAN E. WILLARD, How to Extend the Semantic Tableaux and Cut-Free Versions of the Second Incompleteness Theorem Almost to Robinson’s Arithmetic Q, *Journal of Symbolic Logic* **67**, 465–496 (2002) <http://dx.doi.org/10.2178/jsl/1190150055>
- [14] DAN E. WILLARD, Passive Induction and a Solution to a Paris–Wilkie Open Question, *Annals of Pure and Applied Logic* **146**, 124–149 (2007) <http://dx.doi.org/10.1016/j.apal.2007.01.003>